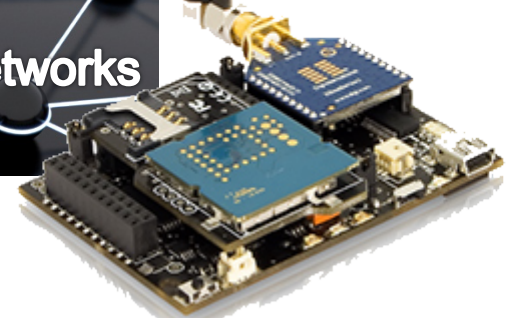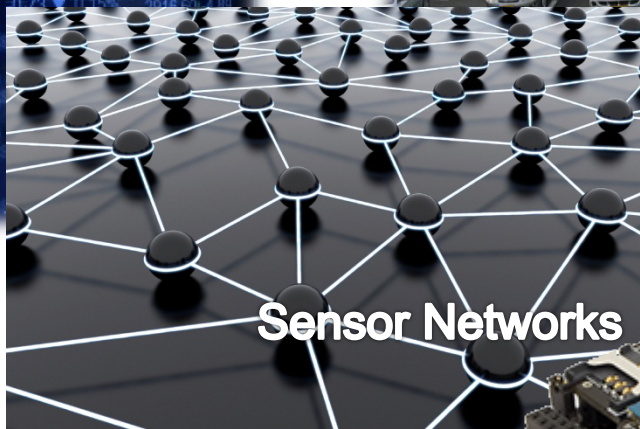# Privacy-Preserving Distributed Stream Monitoring

Arik Friedman[1], Izchak Sharfman[2], Daniel Keren[3], Assaf Schutser[2]

[1] NICTA, Australia    [2] Technion, Israel    [3] Haifa University, Israel

**NICTA**

**NICTA Funding and Supporting Members and Partners**

Australian Government

# Distributed Stream Networks



Financial Data Analysis

Traffic Monitoring Systems

Sensor Networks

## Analyzed data may be personal and sensitive
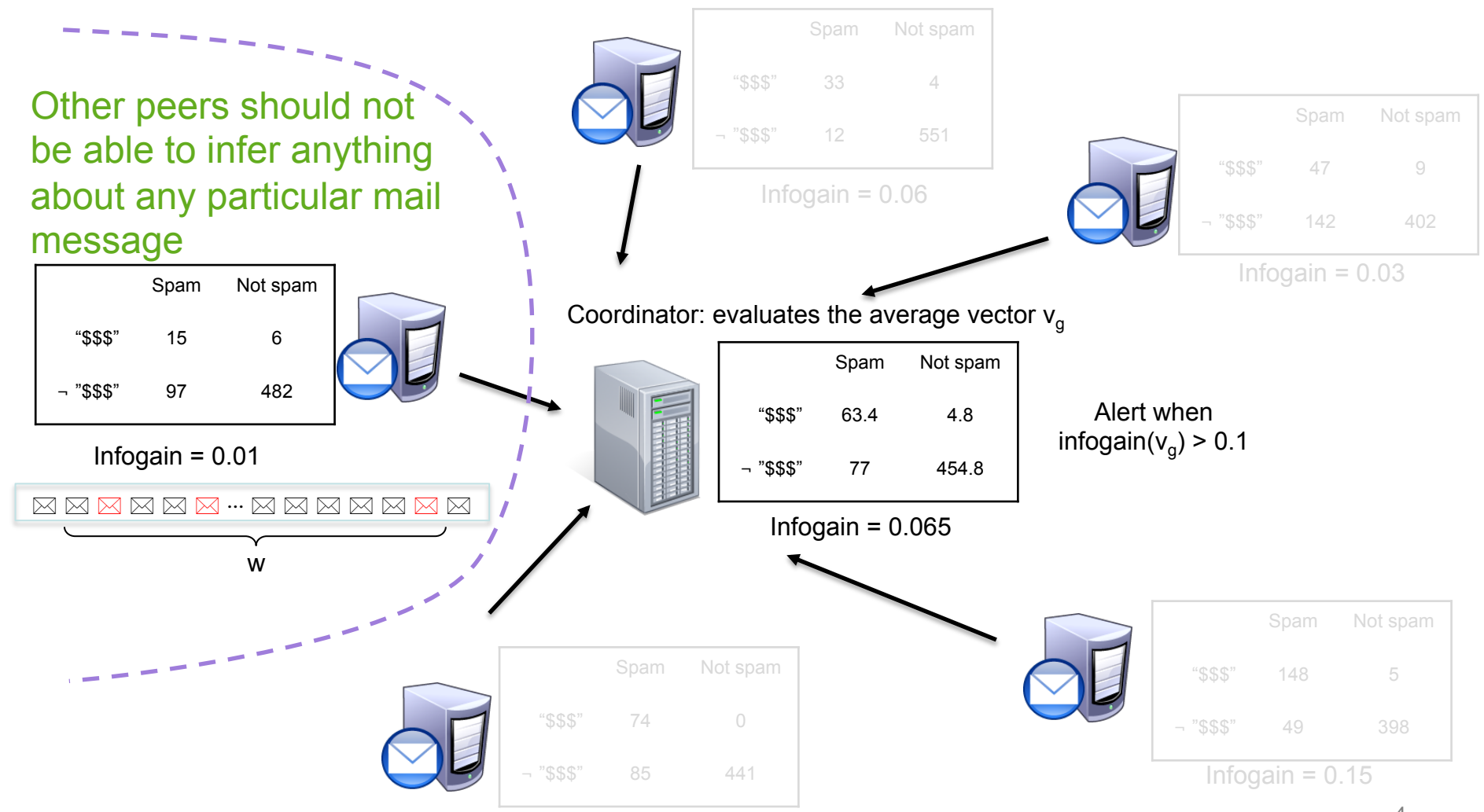
# Related work…

- Continuous monitoring in centralized settings
  - Differential privacy under continual observation [DPNR10]
  - Statistics on sketches [MMNW11]
  - Adaptive sampling [FX12]
- Computation in Distributed settings
  - Distributed noise generation [DKMMN06, CRFG12]
  - Distributed heavy hitters [HKR12]
- Distributed time series data
  - Historical time-series data [RN10]
  - Cryptographic protocols [SCRCS11]
  - Heavy hitters over a sliding window [CLSX12]

This work:

Monitoring complex functions
over statistics derived from streams

# Problem Setting

Other peers should not be able to infer anything about any particular mail message

|  | Spam | Not spam |
|---|---|---|
| "$$$" | 15 | 6 |
| ¬ "$$$" | 97 | 482 |

Infogain = 0.01

⊠⊠⊠⊠⊠⊠ ⋯ ⊠⊠⊠⊠⊠⊠⊠

w

|  | Spam | Not spam |
|---|---|---|
| "$$$" | 33 | 4 |
| ¬ "$$$" | 12 | 551 |

Infogain = 0.06

|  | Spam | Not spam |
|---|---|---|
| "$$$" | 47 | 9 |
| ¬ "$$$" | 142 | 402 |

Infogain = 0.03

Coordinator: evaluates the average vector $v_g$

|  | Spam | Not spam |
|---|---|---|
| "$$$" | 63.4 | 4.8 |
| ¬ "$$$" | 77 | 454.8 |

Infogain = 0.065

Alert when infogain($v_g$) > 0.1

|  | Spam | Not spam |
|---|---|---|
| "$$$" | 74 | 0 |
| ¬ "$$$" | 85 | 441 |

Infogain = 0.08

|  | Spam | Not spam |
|---|---|---|
| "$$$" | 148 | 5 |
| ¬ "$$$" | 49 | 398 |

Infogain = 0.15

From imagination to impact

4

# Problem Setting

Other peers should not be able to infer anything about any particular mail message

|        | Spam | Not spam |
|--------|------|----------|
| "$$$"  | 15   | 6        |
| ¬ "$$$"| 97   | 482      |

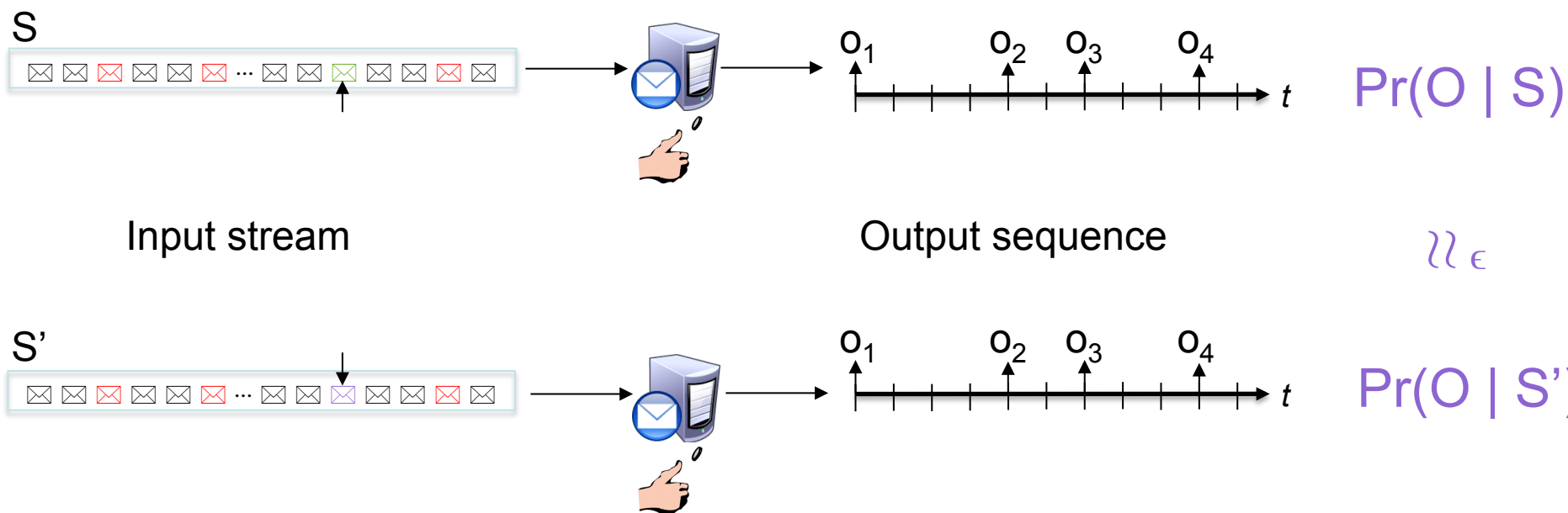Infogain = 0.01

⊠⊠⊠⊠⊠⊠ ... ⊠⊠⊠⊠⊠⊠

w

Cryptographic solutions:

☑ Confidentiality
⊠ Inferences from the output still possibly

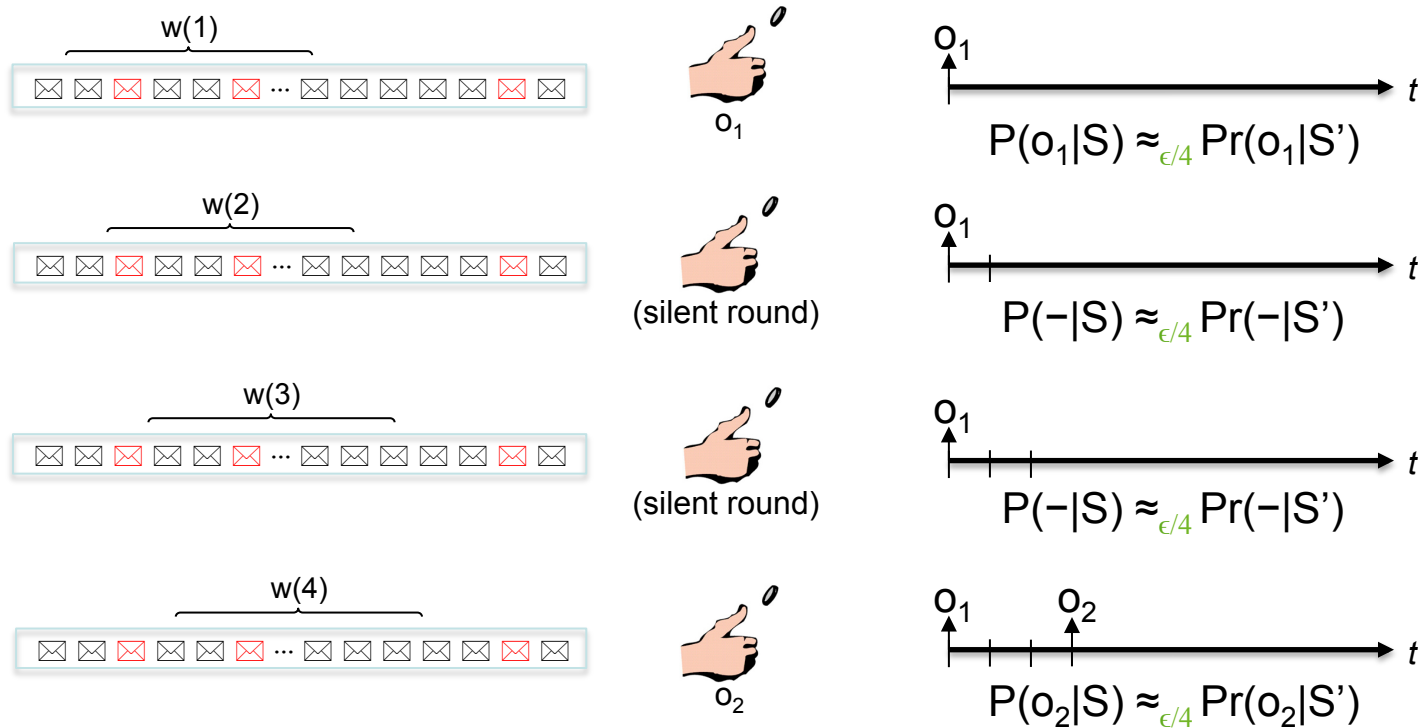⇒Differential privacy addresses such leaks

# Differential privacy [DPNR10]

For any two *adjacent* streams                    and for any output sequence O

S

$o_1$   $o_2$   $o_3$   $o_4$

$t$

$Pr(O \mid S)$

Input stream                                         Output sequence

$\epsilon$

S'

$o_1$   $o_2$   $o_3$   $o_4$

$t$

$Pr(O \mid S')$

Large $\epsilon$ allows bigger difference between the probabilities
$\Rightarrow$ reflects the input more accurately, less private

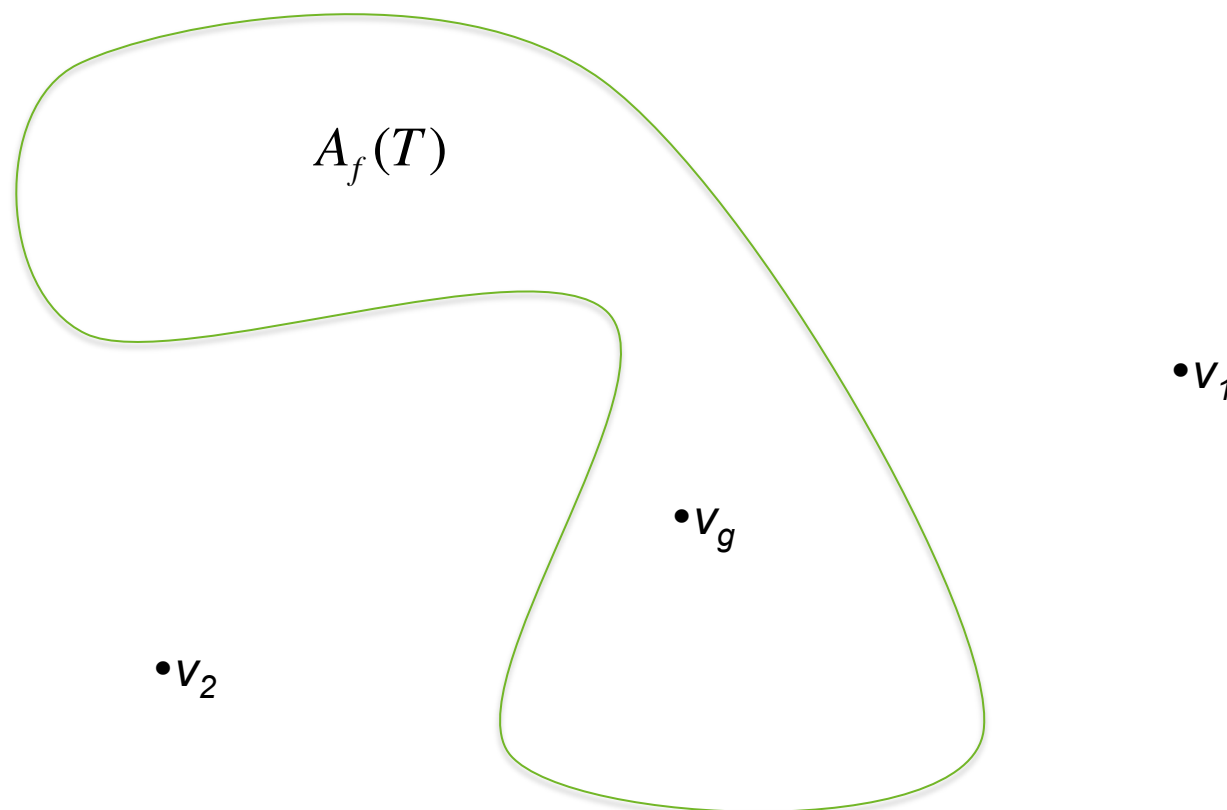# Privacy as a Budget - Naïve Solution



$$\Rightarrow P(o_1 - o_2 | S) \approx_\epsilon Pr(o_1 - o_2 | S')$$

Privacy loss in each time period $\Rightarrow$ wasteful, outputs are not independent
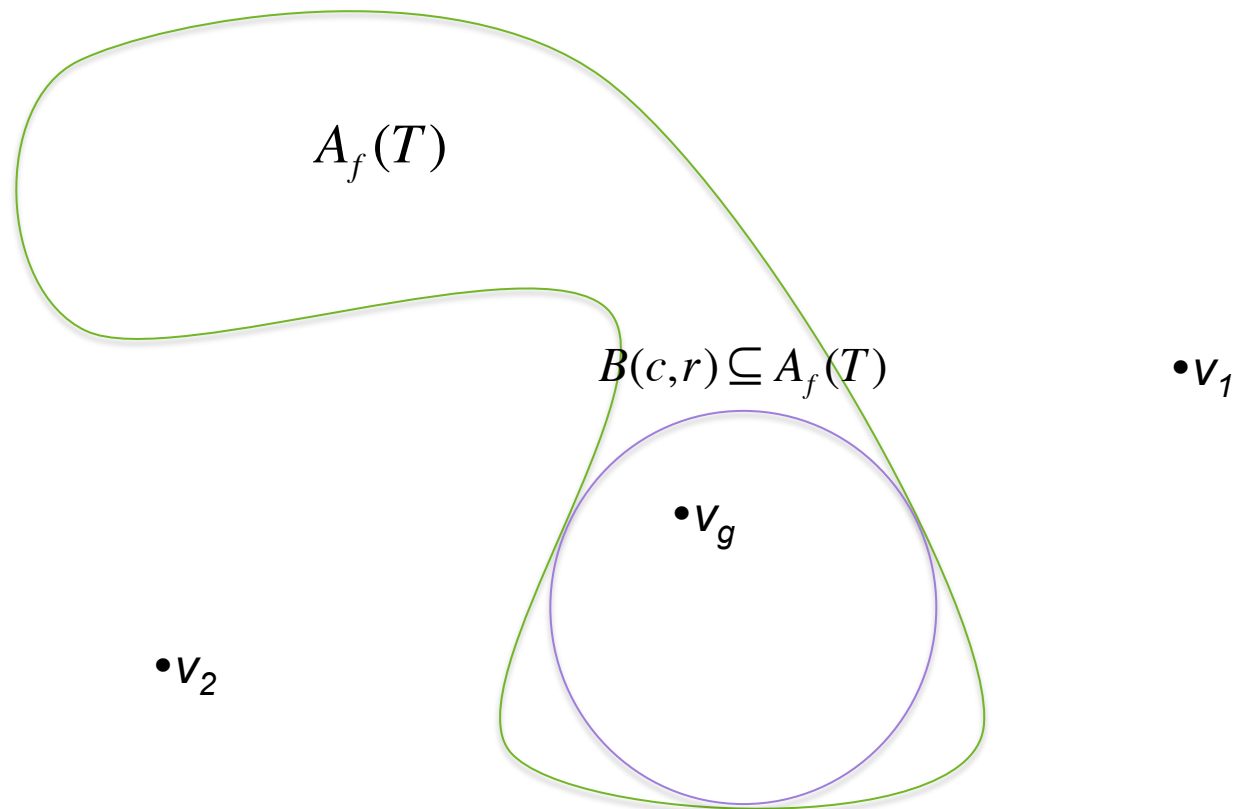Instead, privacy cost can be *amortized*

# Efficient stream monitoring [SSK'06, KSSL'12]

Recall the problem: detect $f(v_g) > T$ for $v_g = \dfrac{1}{k}\sum_k v_i$

The admissible region: $A_f(T) = \left\{ v \mid f(v) \leq T \right\}$

$A_f(T)$

$\bullet v_1$

$\bullet v_g$

$\bullet v_2$

# Efficient stream monitoring [SSK'06, KSSL'12]

Recall the problem: detect $f(v_g) > T$ for $v_g = \dfrac{1}{k}\sum_k v_i$

The admissible region: $A_f(T) = \left\{ v \,\middle|\, f(v) \le T \right\}$

$A_f(T)$

$B(c,r) \subseteq A_f(T)$

$\bullet v_1$

$\bullet v_g$

$\bullet v_2$

# Efficient stream monitoring [SSK'06, KSSL'12]

Recall the problem: detect $f(v_g) > T$ for $v_g = \frac{1}{k}\sum_k v_i$

le region: $A_f(T) = \left\{ v \big| f(v) \leq T \right\}$

Global constraint to local constraints:

$v_g \in B(c,r)$ as long as
$\forall i : v_i \in B(c_i, r)$

$A_f(T)$

$B(c_1, r)$

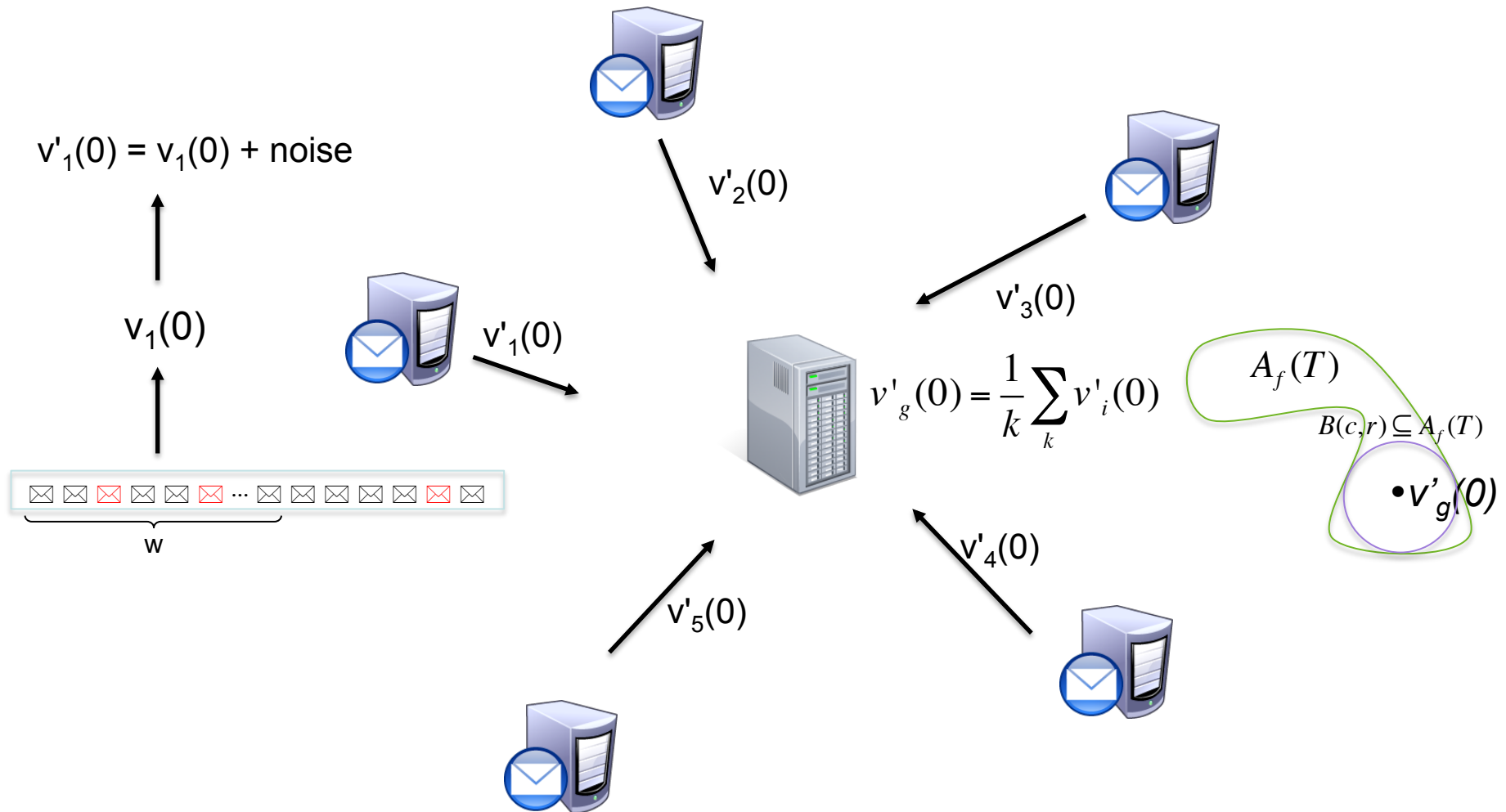$B(c,r) \subseteq A_f(T)$

$\bullet v_1$

$B(c_2, r)$

$\bullet v_g$

$\bullet v_2$

From imagination to impact

$c = \frac{1}{k}\sum_k c_i$

Safe zone
for node 2

Safe zone
for node 1

# Our Algorithm



$v'_1(0) = v_1(0) + noise$

$v_1(0)$

$v'_1(0)$

$v'_2(0)$

$v'_3(0)$

$v'_4(0)$

$v'_5(0)$

w

$$v'_g(0) = \frac{1}{k}\sum_k v'_i(0)$$

$A_f(T)$

$B(c,r) \subseteq A_f(T)$

$\bullet v'_g(0)$

# Our Algorithm

$B(c_1, r')$

$\bullet v_1(0)$

$SZ_1$

$SZ_2$

$SZ_3$

$SZ_4$

$SZ_5$

$$v'_g(0) = \frac{1}{k} \sum_k v'_i(0)$$

$A_f(T)$

$B(c,r) \subseteq A_f(T)$

$\bullet v'_g(0)$

w

From imagination to impact

# Privacy at the Node Level

$B(c_1, r')$

$\bullet v_1(0)$



$r''$

$r'$

$u_1(1)\bullet$  $\bullet v_1(1)$

$c_1$

$u_1(2)\bullet$  $\bullet v_1(2)$

$u_1(3)\bullet$  $\bullet v_1(3)$

$u_1(4)\bullet$  $\bullet v_1(4)$

⊠ ⊠ ⊠ ⊠ ⊠ ⊠ ⋯ ⊠ ⊠ ⊠ ⊠ ⊠ ⊠

w

Evaluating $v_1(t)$ against
the safe zone in Stream S:
t=1: silent round
t=2: silent round
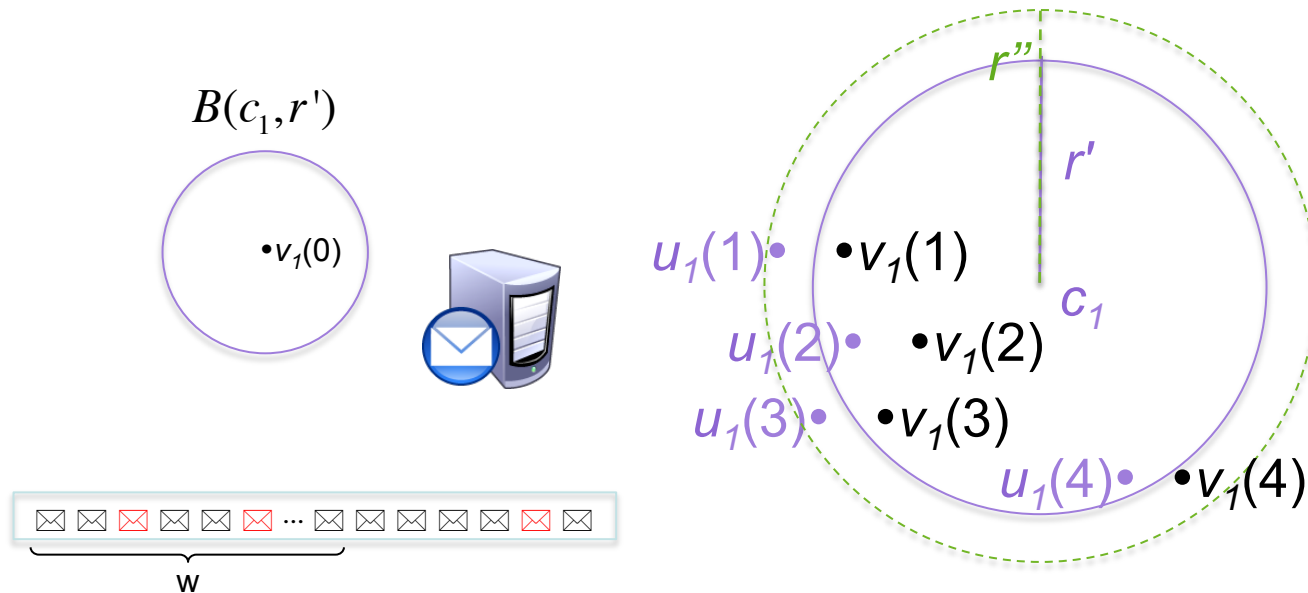t=3: silent round
t=4: safe zone breach

Evaluating $u_1(t)$ against
the safe zone in Stream S':
t=1: ~~silent round~~ breach!

⇒ Addressed by adding
randomness to the safe zone radius (Laplace mechanism)
Pr(silent | S) ≈ Pr(silent | S') because Pr(r') ≈ Pr(r'')

Noise added to the
safe zone will protect
the privacy in all
silent rounds, until a
new safe zone is
assigned!

# Privacy at the Node Level

$B(c_1, r')$

$\bullet v_1(0)$

$r''$

$u_1(1)\bullet$  $\bullet v_1(1)$

$r'$

$u_1(2)\bullet$  $\bullet v_1(2)$  $c_1$

$u_1(3)\bullet$  $\bullet v_1(3)$

$u_1(4)\bullet$  $\bullet v_1(4)$

w

Evaluating $v_1(t)$ against
the safe zone in Stream S:
t=1: silent round
t=2: silent round
t=3: silent round
t=4: safe zone breach

Evaluating $u_1(t)$ against
the safe zone in Stream S':
t=1: silent round
t=2: silent round
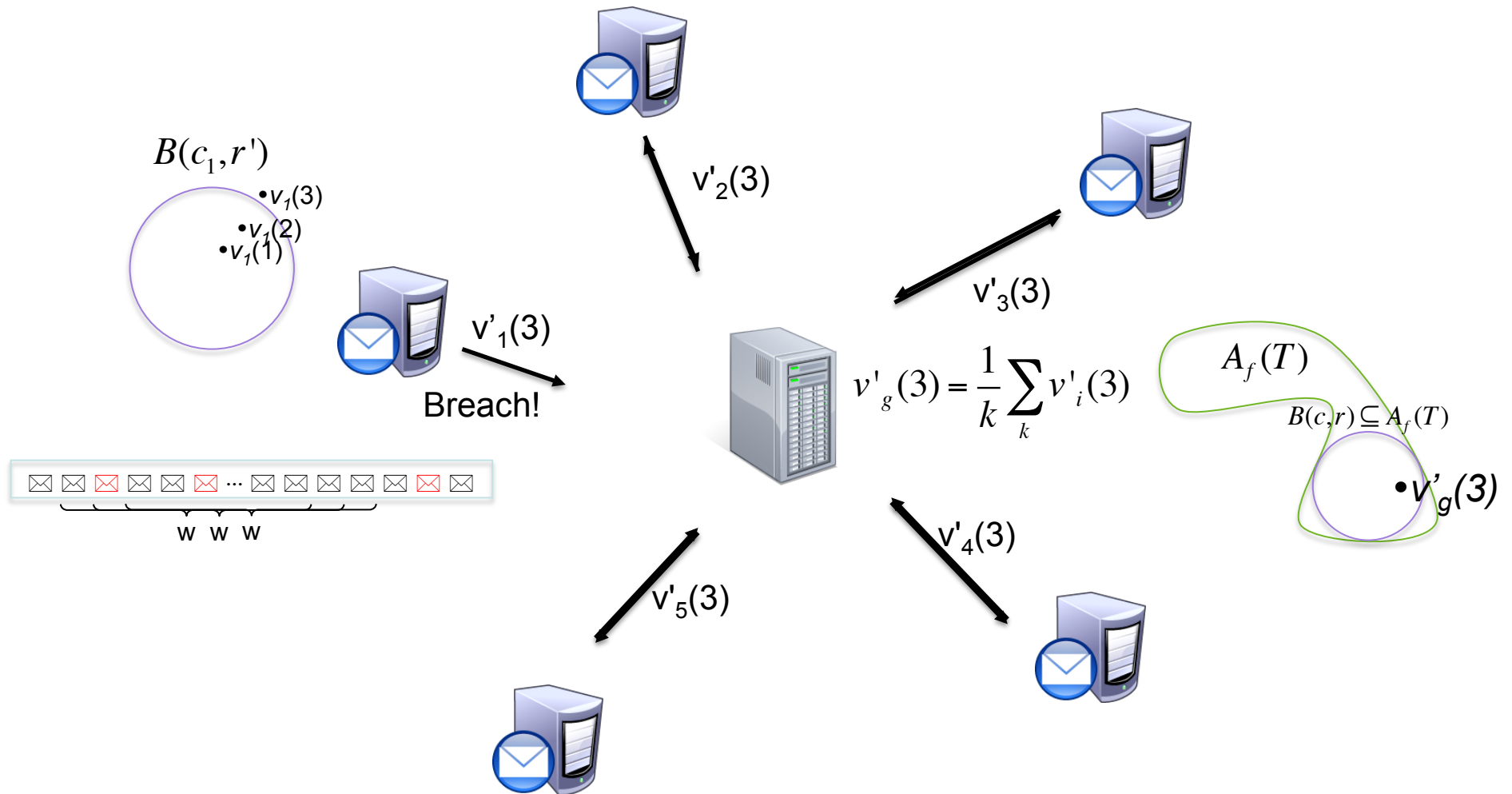t=3: silent round
t=4: ~~safe zone breach~~ silent round
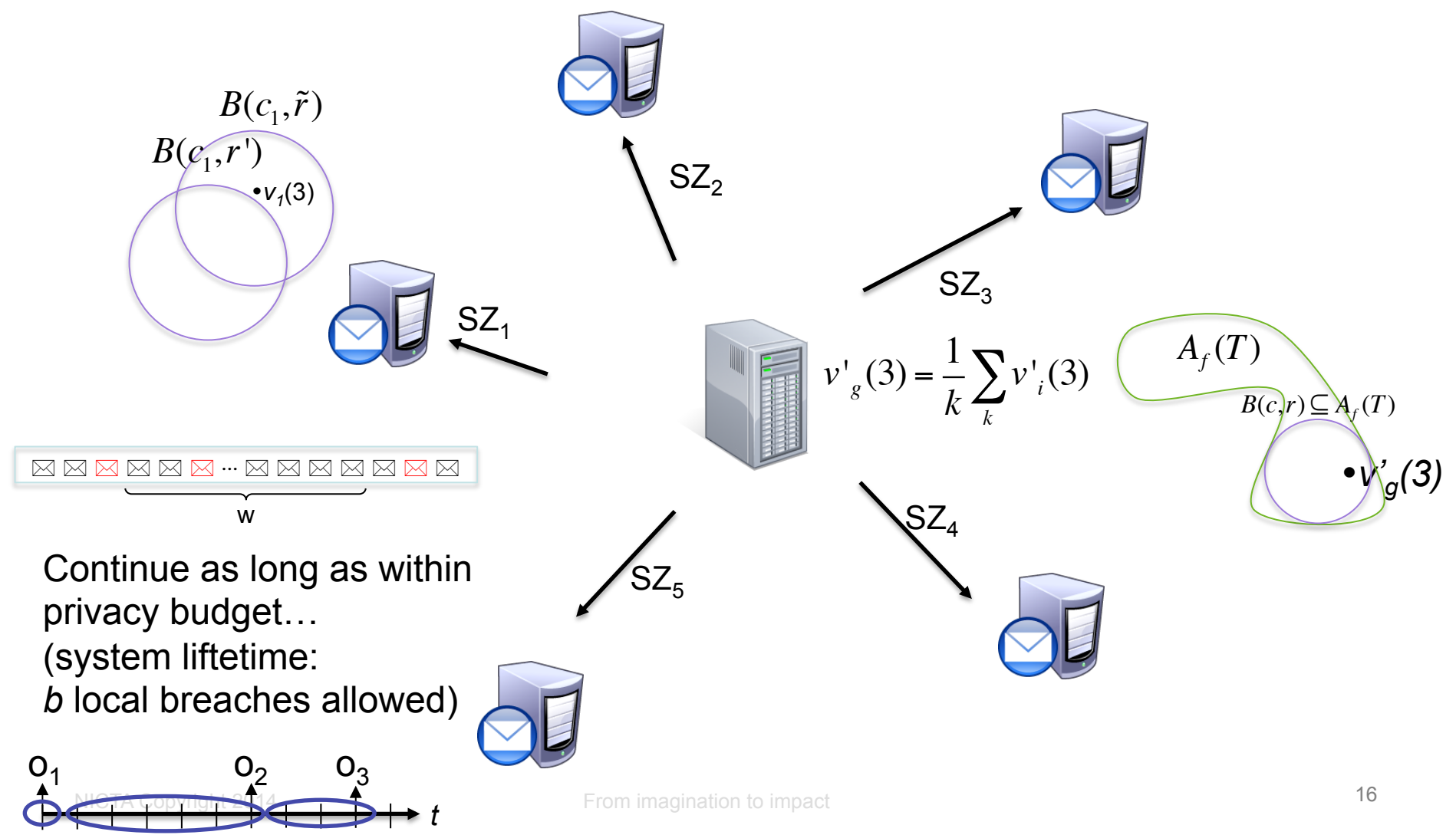
⇒ Addressed by
adding
randomness
(exponential mechanism)
when evaluating

$$v(t) \in_\varepsilon B(c, r')$$

# Our Algorithm



$B(c_1,r')$

$\bullet v_1(3)$
$\bullet v_1(2)$
$\bullet v_1(1)$

$v'_2(3)$

$v'_3(3)$

$v'_1(3)$

Breach!

$v'_g(3) = \dfrac{1}{k}\sum_k v'_i(3)$

$A_f(T)$

$B(c,r) \subseteq A_f(T)$

$\bullet V'_g(3)$

w  w  w

$v'_4(3)$

$v'_5(3)$

# Our Algorithm



$$B(c_1, \tilde{r})$$

$$B(c_1, r')$$

$\bullet v_1(3)$

$SZ_2$

$SZ_3$

$SZ_1$

$$v'_g(3) = \frac{1}{k} \sum_k v'_i(3)$$

$$A_f(T)$$

$$B(c,r) \subseteq A_f(T)$$

$\bullet v'_g(3)$

$SZ_4$

$SZ_5$

w

Continue as long as within
privacy budget...
(system lifetime:
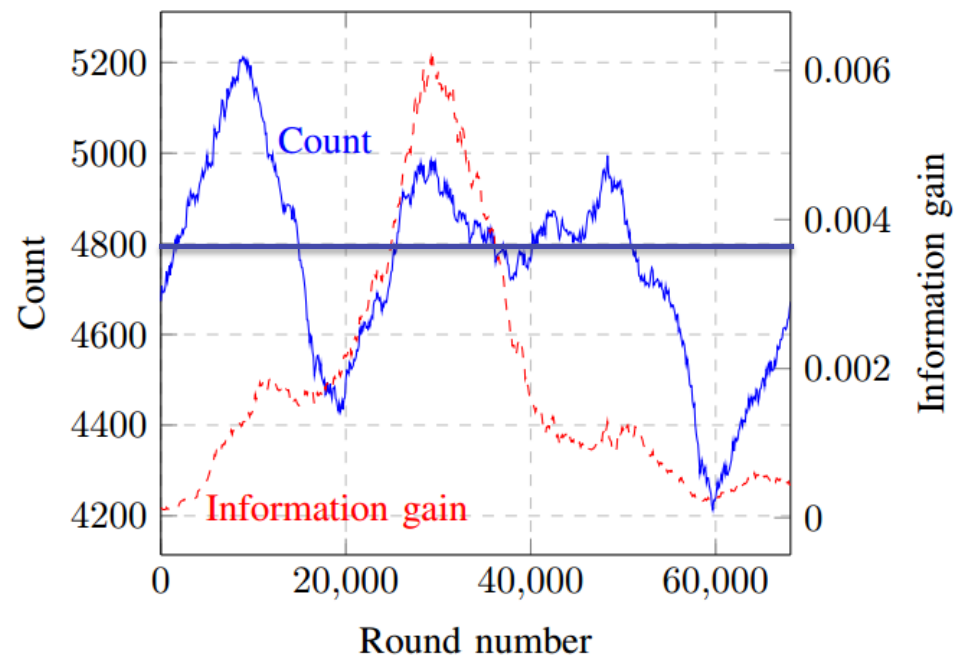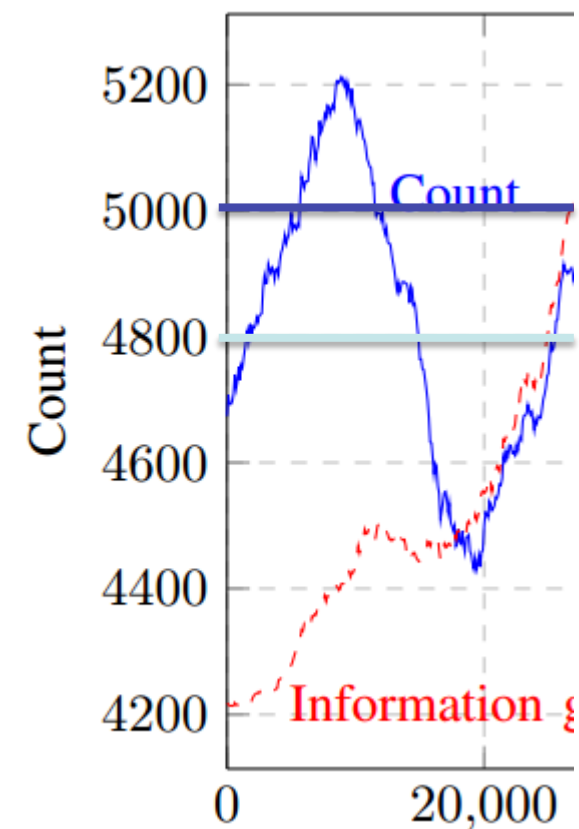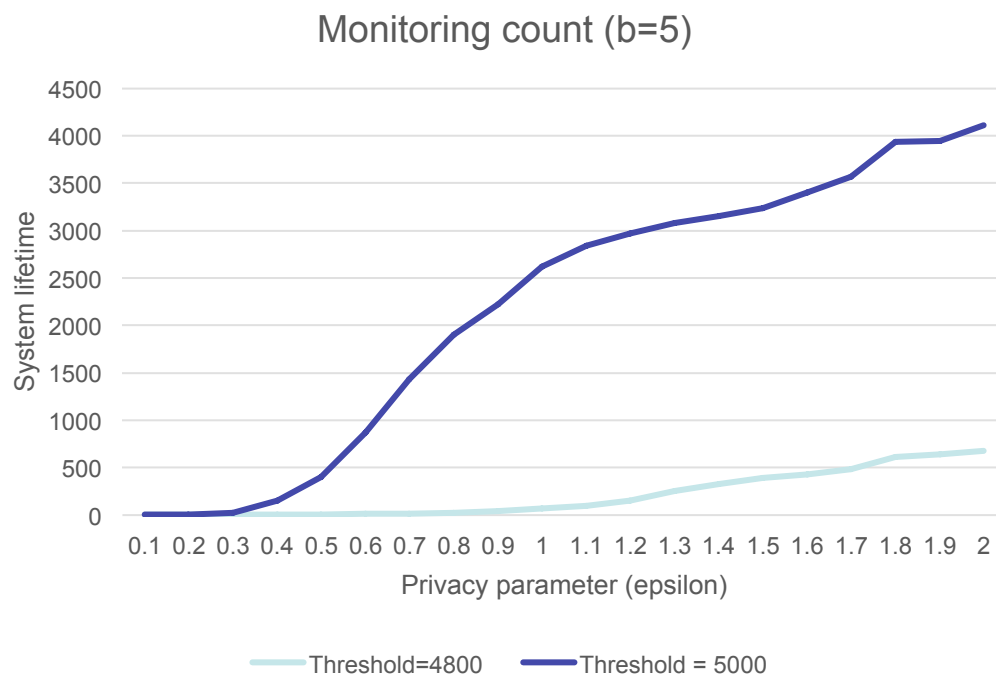$b$ local breaches allowed)

$o_1$    $o_2$    $o_3$

$t$

# Experimental evaluation

Reuters corpus:

- 781,265 labelled news stories
- Distributed by round robin between 10 nodes
- Each node monitors a window of 10,000 stories
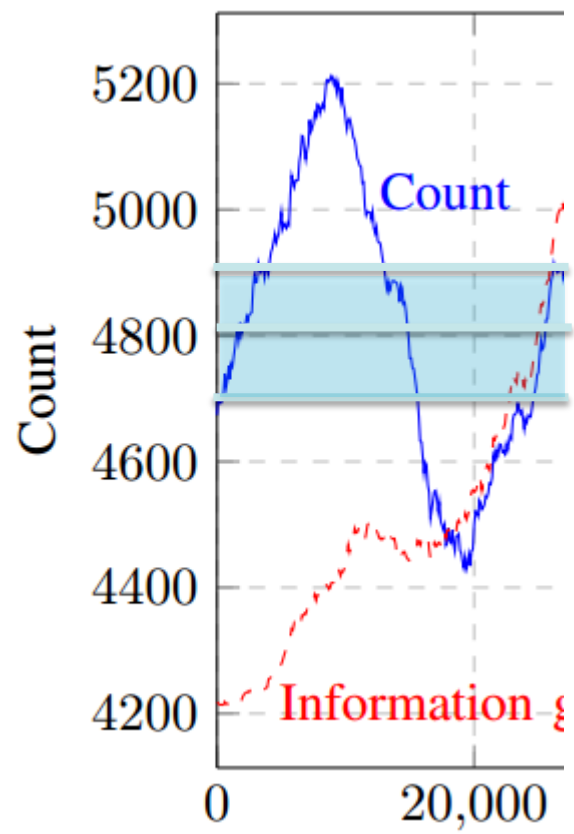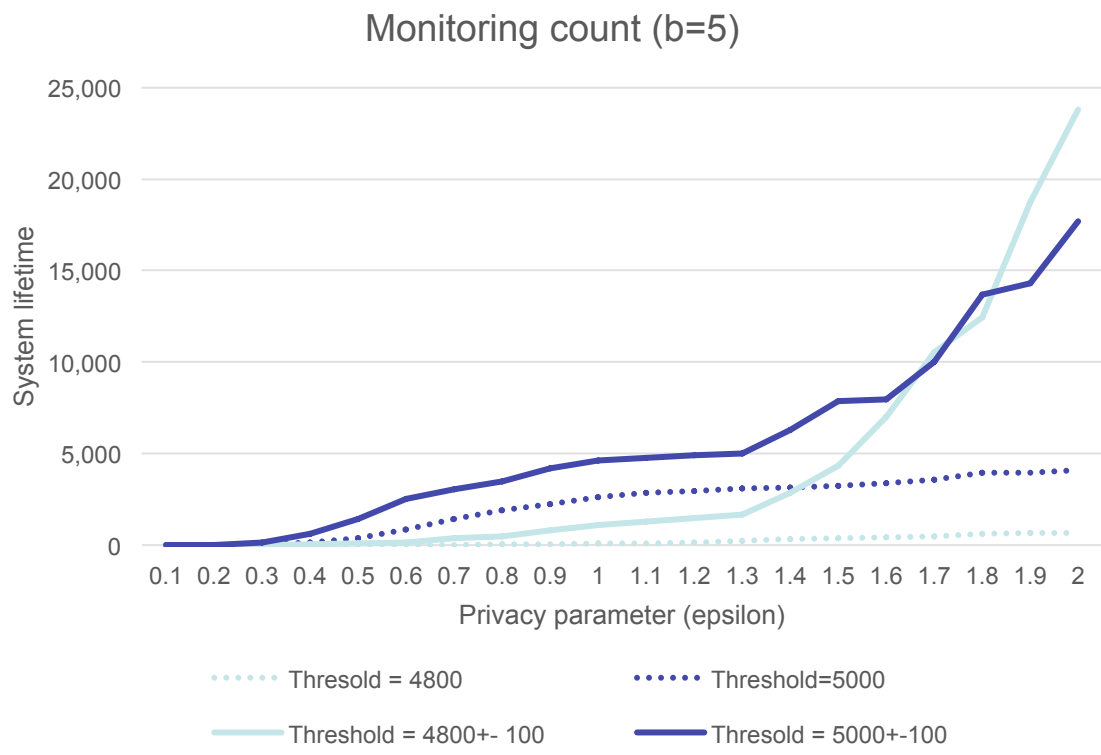- "CCAT" category denotes spam, "febru" feature a monitored term

# Monitoring count



Monitoring count (b=5)



Likelihood of local breach higher
when closer to the threshold

# Adding error margins



Monitoring count (b=5)



Error margins trade accuracy for longer system lifetime

# Additional results in the paper…

- Infogain evaluation
  - Tradeoff between System lifetime, threshold and privacy: we pay for privacy mainly when close to the threshold.

- Error margins trade-offs

- Violation rounds (local breaches $b$) trade-off

- Costs of distributed vs. centralized

From imagination to impact

# Summary and future directions

NICTA

Communication efficiency translates
to better privacy

- Possible enhancements:
  - Local communication between nodes could allow further mitigation of privacy loss
  - Prediction models that tailor safe zones to nodes can reduce the probability of local breaches
  - As the processing window advances, the privacy budget can be replenished

From imagination to impact

# Thank you

NICTA @ Sydney
(we hire!)

From imagination to impact